

# Détection des cultures et des mauvaises herbes

ELHASSOUNI Anas, KINDA Abdoul Latif, MAHARAVO Tefinjanahary Anicet  
Superviseur : Prof. Anass BELCAID

10 janvier 2025

## Résumé

L'agriculture de précision et en particulier l'application de l'intervention automatisée contre les mauvaises herbes représentent un domaine de recherche de plus en plus essentiel, car les considérations de durabilité et d'efficacité deviennent de plus en plus pertinentes. Alors que les potentiels des réseaux neuronaux convolutionnels pour les tâches de détection, de classification et de segmentation ont été démontrés avec succès dans d'autres domaines d'application, ce domaine relativement nouveau ne dispose actuellement pas de la quantité et de la qualité requises de données d'entraînement pour une approche aussi axée sur les données. Par conséquent, nous proposons un nouvel ensemble de données d'images à grande échelle spécialisé dans l'identification fine de 74 espèces de cultures et de mauvaises herbes pertinentes, en mettant l'accent sur la variabilité des données. Nous fournissons des annotations de boîtes englobantes étiquetées, de masques sémantiques et de positions de tiges pour environ 112 000 instances dans plus de 8 000 images haute résolution de sites agricoles réels et de parcelles extérieures spécifiquement cultivées de types de mauvaises herbes rares. De plus, chaque échantillon est enrichi d'un ensemble complet de méta-annotations concernant les conditions environnementales et les paramètres d'enregistrement. Nous menons en outre des expériences de référence pour plusieurs tâches d'apprentissage sur différentes variantes de l'ensemble de données afin de démontrer sa polyvalence et de fournir des exemples de schémas de cartographie utiles pour adapter les données annotées aux exigences d'applications spécifiques. Au cours de l'évaluation, nous démontrons en outre comment l'intégration de plusieurs espèces de mauvaises herbes dans le processus d'apprentissage augmente la précision de la détection des cultures. Dans l'ensemble, l'évaluation démontre clairement que notre ensemble de données représente une étape essentielle pour combler le manque de données et promouvoir de nouvelles recherches dans le domaine de l'agriculture de précision.

## 1 Introduction

En période de croissance démographique mondiale, l'industrie agricole est confrontée à une demande croissante de cultures vivrières, tandis que la réduction de son impact sur l'environnement et la santé humaine devient simultanément une priorité. L'une des tâches offrant le plus grand potentiel est l'élimination efficace des mauvaises herbes, qui concurrencent les cultures cultivées pour des ressources telles que la lumière du soleil, l'eau, l'espace et les nutriments, réduisant ainsi considérablement leur rendement global. Au lieu de la technique couramment utilisée consistant à appliquer des herbicides à grande échelle sur l'ensemble de la surface cultivée, l'agriculture de précision, en particulier sous la forme de localisation et de classification automatiques des cultures et des

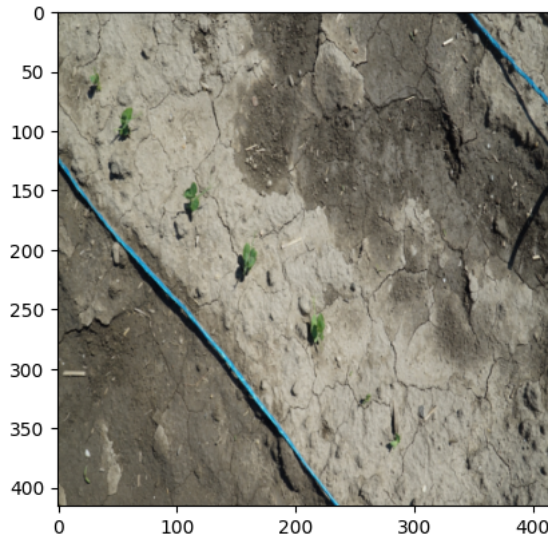


FIGURE 1 – Image du dataset

mauvaises herbes, offre la possibilité d’augmenter considérablement l’efficacité et la durabilité du processus en appliquant avec précision des herbicides à des plantes spécifiques ou même en éliminant les mauvaises herbes mécaniquement sans l’utilisation de produits chimiques. En outre, des méthodes similaires peuvent être utilisées pour d’autres tâches telles que le suivi de la croissance, la récolte automatique et même la reconnaissance des ravageurs ou des maladies des plantes. Les réseaux neuronaux convolutionnels (CNN) et les Vision Transformers], qui ont récemment émergé, présentent des approches très prometteuses pour les tâches d’agriculture de précision telles que la segmentation, la classification fine ainsi que la détection des plantes et des tiges. Cependant, leur application robuste nécessite une quantité importante de données d’apprentissage annotées très variables, qui ne sont pas encore disponibles dans ce domaine de recherche. De plus, peu d’importance est accordée à la pertinence des combinaisons de plusieurs cultures et en particulier de mauvaises herbes, qui, à notre avis, devraient être exploitées pour augmenter les performances de détection dans des scénarios réels. Pour répondre à ce besoin de recherche, notre travail propose les contributions suivantes :

- Nous fournissons un nouvel ensemble de données à grande échelle pour l’agriculture de précision, composé d’images réelles très variables et d’annotations multimodales pour un ensemble riche de catégories de cultures et de mauvaises herbes.
- Nous démontrons sa polyvalence en entraînant et en évaluant plusieurs tâches d’apprentissage, notamment la détection d’objets.

## 2 Domaines d’applications

De nombreuses équipes de partout dans le monde travaillent ou ont déjà travaillé sur des projets similaires. L’étude de leurs travaux est une mine de connaissances utiles à notre recherche. Pour la deuxième partie de notre projet, voici quelques exemples d’applications concrètes et d’équipes de recherche qui ont travaillé sur la détection des cultures et des mauvaises herbes en utilisant des réseaux de neurones convolutifs (CNN) et d’autres approches d’intelligence artificielle

## **2.1 Université de Bonn, Allemagne**

L'université de Bonn a mené des recherches approfondies sur la détection des cultures et des mauvaises herbes dans le cadre de l'agriculture de précision. Leur équipe composée de Andres Milioto, Philipp Lottes, Cyrill Stachniss a travaillé sur le développement de modèles de CNN et de Vision Transformers pour la segmentation et la classification des plantes en conditions réelles. Ils ont particulièrement travaillé sur la détection de carottes et des mauvaises herbes. Ils ont utilisé des techniques d'annotation multimodale pour mieux comprendre les caractéristiques des plantes et des mauvaises herbes dans différents environnements agricoles. Ces travaux démontrent comment l'IA peut être utilisée pour optimiser les traitements aux herbicides de manière ciblée.

## **2.2 Université de l'Illinois, États-Unis**

L'université de l'Illinois a exploré l'utilisation de CNN pour la détection et la classification des mauvaises herbes dans les cultures de maïs et de soja. Leur recherche se concentre sur des systèmes de détection embarqués pour les robots agricoles, capables d'identifier et de traiter automatiquement les mauvaises herbes, réduisant ainsi la nécessité de recourir aux produits chimiques. Ils ont aussi développé un jeu de données spécifique pour les cultures américaines, qui peut servir comme base pour améliorer la précision des modèles dans différents environnements.

## **2.3 Institut de Recherche pour le Développement (IRD), France**

En collaboration avec des équipes agricoles locales, l'IRD utilise des modèles CNN pour analyser les images de champs cultivés dans des régions tropicales, particulièrement en Afrique de l'Ouest. Leurs travaux se concentrent sur l'analyse multispectrale et la détection des mauvaises herbes dans les cultures de riz et de mil. L'objectif est de créer des outils abordables pour les agriculteurs, permettant de surveiller les champs et de réduire les pertes de rendement dues aux mauvaises herbes.

Ces exemples montrent comment différentes approches et technologies (CNN, Vision Transformers, drones, robots) sont utilisées pour optimiser les processus agricoles, réduire la dépendance aux herbicides et promouvoir des pratiques plus durables.

# **3 Le dataset CropOrWeed**

Notre ensemble de données proposé constitue une étape essentielle pour combler le manque de données dans l'agriculture de précision. Les données d'image ainsi que les annotations multimodales sont disponibles pour la recherche universitaire et sont destinées à être améliorées en collaboration avec la communauté pour augmenter progressivement la diversité en ajoutant des échantillons collectés dans le monde entier. Les données ont été rassemblées par AIT Austrian Institute of Technology qui a fait un immense travail dans la collecte ainsi que les annotations.

## **3.1 Conception du dataset**

Les données ont été recueillies à l'aide d'un appareil photo reflex numérique semi-professionnel équipé d'un capteur plein format. Ce choix d'équipement a permis de répondre à deux exigences : la qualité d'image et la mobilité nécessaire pour les prises de vue sur le terrain. Toutes les images sont capturées manuellement en mode d'exposition



	Exp	App	Total
<b>Sessions</b>			
Enregistrées	1 363	738	2 101
Annotées	665	264	929
<b>Images</b>			
Enregistrées	22 597	21 217	43 814
Annotées	4 990	3 044	8 034
<b>Instances</b>			
Annotées	66 877	45 076	111 953

TABLE 1 – Résumé des données.

trouvent dans des fichiers CSV du même nom que les images permettant les correspondances. Chaque image contient zéro ou plusieurs plantes ; chaque fichier csv contient au moins une ligne avec sept colonnes.

- Les deux premières représentant les coordonnées du coin supérieur gauche du rectangle entourant l'image
  - Les deux suivantes représentant les coordonnées du coin inférieur droit du rectangle entourant l'image
  - La cinquième colonne représente la classe de l'objet. Le sol ou tout autre objet n'étant pas une plante ou une mauvaise herbe a la classe 255.
  - Les deux dernières colonnes étant les coordonnées du centre du rectangle.
- Chaque ligne représente un objet se trouvant dans l'image



FIGURE 3 – Image du dataset

### 3.3 Division du dataset pour l'apprentissage

Les variantes de l'ensemble de données définies sont utilisées pour former et comparer des modèles spécialisés pour plusieurs tâches d'apprentissage et scénarios d'application. Pour l'entraînement avec YoLoV3, la répartition générée aléatoirement entre les données d'entraînement et de test de 80 :20 est appliquée pour garantir l'indépendance même lors de la combinaison de tâches d'apprentissage et l'entraînement avec Fast-RCNN a été fait avec une division de 70 :30.

L'entraînement est effectué sur un système avec un GPU NVIDIA TESLA P100 fourni par Kaggle.

Tout au long de l'entraînement, nous avons utilisé ADAM pour l'optimisation et nous avons opté pour Pytorch.

Nous avons ainsi utilisé les modèles YoLoV3 et Faster-RCNN qui ont été à partir de zéro, aucun modèle pré-entraîné n'a ainsi été utilisé.

## 4 Les modèles utilisés

### 4.1 YoLoV3

YOLOv3 est un modèle de détection d'objets en temps réel qui améliore les versions précédentes de YOLO (YOLOv1 et YOLOv2). Développé par Joseph Redmon et al. en 2018, YOLOv3 est connu pour son équilibre entre vitesse et précision dans des tâches complexes de détection d'objets.

- **Détection en temps réel** : YOLOv3 divise une image en une grille et effectue simultanément la localisation des objets (boîtes englobantes) et leur classification. Cela le rend extrêmement rapide pour des applications en temps réel, comme la vidéosurveillance et la conduite autonome.

- **Multi-échelle (Feature Pyramid)** : YOLOv3 détecte les objets à trois échelles différentes. Chaque échelle utilise des cartes de caractéristiques extraites à différents niveaux de profondeur du réseau, ce qui améliore la détection d'objets de différentes tailles ; les petits objets sont détectés dans des cartes de caractéristiques fines tandis que les objets moyens et grands sont détectés dans des cartes de caractéristiques plus profondes.

- **Architecture basée sur Darknet-53** : YOLOv3 repose sur Darknet-53, un réseau de neurones convolutifs (CNN) soit 53 couches convolutionnelles et est optimisé grâce à des connexions résiduelles (ResNet) pour éviter le problème de dégradation dans les réseaux profonds.

- **Sortie du réseau** : À chaque échelle, YOLOv3 génère des prédictions sous forme de tenseurs avec les informations suivantes pour chaque boîte englobante, coordonnées (x, y, w, h) pour la position du centre et la taille de la boîte, un score de confiance qui est la probabilité qu'une boîte contienne un objet et des scores des classes qui sont les probabilités que l'objet appartienne à une certaine catégorie (par exemple, mais, betterave).

- **Ancres (Anchors)** : YoloV3 utilise des boîtes d'ancrage pour prédire efficacement les objets de différentes tailles et formes. Ces boîtes sont adaptées aux distributions des objets dans l'ensemble d'entraînement.

- **Fonction de perte** : YOLOv3 utilise une fonction de perte qui combine le calcul de l'erreur des coordonnées pour améliorer la localisation des boîtes, l'erreur de confiance pour différencier les boîtes contenant des objets des fausses prédictions, l'erreur de classification pour associer la boîte à la bonne classe.

### 4.2 Faster-RCNN

Faster-R-CNN est une amélioration de R-CNN (Region-based Convolutional Neural Network) et Fast-RCNN, proposée par Ross Girshick en 2015. Ce modèle est conçu pour accélérer le processus de détection d'objets tout en maintenant une précision élevée. C'est un modèle largement utilisé pour la détection d'objets et la segmentation d'images dans



des applications de vision par ordinateur.

- **Réseaux de Proposition de Régions (RPN)** : génère des propositions de régions rectangulaires avec des scores de probabilité d'appartenance à un objet à partir d'une image. Ce processus utilise un réseau entièrement convolutionnel pour partager les calculs avec un détecteur Fast R-CNN, basé sur des couches convolutionnelles, comme celles du modèle VGG-16. Un petit réseau glissant sur la carte des caractéristiques produit des propositions en prenant des fenêtres de taille  $n \times n$  (avec  $n = 3$ ), qui sont réduites à des vecteurs de 512 dimensions (VGG) avec une activation ReLU. Ces vecteurs alimentent deux couches parallèles pour la régression des boîtes (reg) et leur classification (cls). L'approche glissante permet de partager les couches entièrement connectées sur toutes les positions spatiales, ce qui est efficacement implémenté avec des couches convolutionnelles  $n \times n$  et  $1 \times 1$ .

- **Les ancres** : À chaque position de la fenêtre glissante, plusieurs propositions de régions sont simultanément prédites, avec un maximum de  $k$  propositions. La couche de régression génère  $4k$  sorties pour encoder les coordonnées des  $k$  boîtes, tandis que la couche de classification produit  $2k$  scores pour estimer la probabilité d'appartenance à un objet. Ces  $k$  propositions sont définies par rapport à  $k$  ancres, des boîtes de référence centrées sur la fenêtre glissante et associées à différentes échelles et rapports d'aspect. Avec 3 échelles et 3 rapports d'aspect par défaut, cela donne  $k = 9$  ancres par position. Pour une carte de caractéristiques de taille  $W \times H$ , cela génère  $W \times H \times \text{ancres}$ .

La méthode des ancres offre une solution efficace pour gérer plusieurs échelles et rapports d'aspect. Contrairement aux pyramides d'images ou de caractéristiques, coûteuses en temps, et aux fenêtres glissantes de tailles variées, notre approche utilise une "pyramide d'ancres" avec des cartes de caractéristiques et des images à une seule échelle. Cette méthode, basée sur des filtres de taille unique, est plus rentable tout en conservant une grande efficacité.

- **Fonction de perte** : Pour l'entraînement des RPN, chaque ancrage reçoit une étiquette binaire indiquant s'il représente un objet ou non. Les étiquettes positives sont attribuées à deux types d'ancrages : (i) ceux ayant le plus grand recouvrement Intersection-over-Union (IoU) avec une boîte de vérité de terrain, ou (ii) ceux ayant un IoU supérieur à 0,7 avec une boîte de vérité de terrain. Une boîte de vérité de terrain peut attribuer des étiquettes positives à plusieurs ancrages. Bien que la deuxième condition soit généralement suffisante, la première est conservée pour éviter les cas où aucun échantillon positif ne serait trouvé.

Les ancres sont considérés comme négatifs si leur IoU est inférieur à 0,3 pour toutes les boîtes de vérité de terrain. Les ancres qui ne sont ni positifs ni négatifs sont exclus de l'objectif d'entraînement.

Avec ces définitions, une fonction de perte multitâche, similaire à celle de Fast R-CNN, est utilisée. Elle combine la classification et la régression pour optimiser l'entraînement des propositions de régions.

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*). \quad (1)$$

FIGURE 4 – Calcul de la perte

## 5 Les résultats de l'entraînement des modèles

### 5.1 YOLOV3

Pour le modèle YOLOV3, les données ont d'abord été préparées pour correspondre aux données d'entrées de YOLO. C'est ainsi que nous avons créé de nouveaux fichiers au format .txt contenant les coordonnées du centre de chaque boîte, sa longueur et largeur ainsi que la classe de l'objet.

Nous nous sommes retrouvés avec 7800 images au format .jpg et fichiers d'annotations .txt. Pour l'entraînement, 6427 images aléatoires et leurs annotations ont été utilisées correspondant à la division 80 :20.

Une stratégie d'échantillonnage est suivie, où chaque mini-lot est constitué d'un ensemble d'images contenant des objets annotés. YOLOv3 divise l'image en une grille et attribue les responsabilités de détection aux cellules en fonction de l'appartenance des objets aux ancres définies.

Pour éviter un biais en faveur des échantillons négatifs (les cellules qui ne contiennent pas d'objets), YOLOv3 ajuste la pondération des pertes en fonction des proportions positives et négatives dans l'image. Les échantillons positifs sont déterminés par les ancres ayant le meilleur recouvrement Intersection-over-Union (IoU) avec les boîtes de vérité de terrain.

Les poids des couches du réseau sont initialisés aléatoirement en utilisant une distribution gaussienne centrée sur zéro, avec un écart-type typique pour stabiliser l'entraînement. Le taux d'apprentissage suit une politique de décroissance par paliers : un taux initial élevé (par exemple, 0,0001) est appliqué pendant plusieurs itérations, puis réduit à 0,0001 à des étapes définies. Le moment est fixé à 0,9 et une régularisation par décroissance de poids de 0,0005 est utilisée pour prévenir le surapprentissage.

YOLOv3 optimise une fonction de perte composée de trois termes : la régression des coordonnées des boîtes englobantes, la prédiction des classes, et la confiance associée aux propositions. Cette approche unifiée permet une détection rapide et précise dans un cadre de réseau unique. Au bout de l'entraînement sur quarante epochs, nous avons obtenu une perte totale c'est-à-dire la somme de la perte de classification, de confiance et des boxes de 13.4 et une précision de 23%. Les modèles YOLOv3 et YOLOv5s ont été entraînés par

YOLOv3	YOLOv5s	YOLOv5l
23.0	47.4	71.5

TABLE 2 – Précision mAP des différents modèles.

notre équipe. La précision moyenne établie a été calculée en se reposant sur des mesures



d'intersection\_over\_union en omettant les objets qui ne sont pas des plantes ou l'arrière-plan.

## 5.2 Faster-RCNN

- **Entraînement des RPN** : Le RPN est entraîné par rétropropagation et descente de gradient stochastique (SGD) avec une stratégie d'échantillonnage centrée sur l'image. Chaque mini-lot contient 256 ancres, avec un ratio maximum de 1 :1 entre exemples positifs et négatifs. Les poids des couches sont initialisés aléatoirement avec une distribution gaussienne (écart-type de 0,01), un taux d'apprentissage de 0,001 pendant 60 000 mini-lots, puis réduit à 0,0001 pour les 20 000 suivants. Le moment est fixé à 0,9 et la décroissance de poids à 0,0005.

- **Implémentation** : Les images sont redimensionnées pour que leur côté le plus court mesure 600 pixels. Les ancres utilisent 3 échelles ( $128^2$ ,  $256^2$ ,  $512^2$  pixels) et 3 rapports d'aspect (1 :1, 1 :2, 2 :1). Pendant l'entraînement, les ancres dépassant les limites de l'image sont ignorées pour éviter des erreurs importantes. Une image typique contient environ 6000 ancres utilisables après filtrage.

Pendant les tests, des boîtes de propositions dépassant les limites sont ajustées pour s'adapter aux bords de l'image. Une suppression des non-maximaux (NMS) est appliquée avec un seuil IoU de 0,7, réduisant les propositions à environ 2000 par image, sans affecter la précision de détection finale. Faster R-CNN est entraîné avec 2000 propositions, mais le nombre peut varier lors des tests.

Au bout de l'entraînement, nous obtenons une perte de classification RPN de 0,1951, une perte localisation RPN de 0,4494, une perte de classification Faster-RCNN de 0,2964, une perte localisation Faster-RCNN de 5,9842.

Nous avons obtenu une précision moyenne établie à 46%.

## 6 Discussion

Au bout de l'entraînement de chaque modèle, il en ressort que le meilleur modèle entre les deux est le *Faster — RCNN* parmi les modèles non-préentraînés. La prochaine étape pour nous sera d'utiliser des modèles plus performants comme YOLOV9 pré-entraînés pour profiter de la modernité de ce modèle. Une étape supérieure serait également d'augmenter la taille du dataset en y ajoutant davantage d'images annotées et d'utiliser des stratégies de mise en échelle. Une autre étape sera également d'augmenter le nombre de classe afin de pouvoir détecter et classer davantage de plantes.

L'émergence de tels modèles de détection permettra dans le futur de donner un nouvel élan à l'agriculture de précision en favorisant la croissance des plantes et améliorer les récoltes.

Au regard de la croissance démographique continue, l'agriculture de précision va prendre de plus en plus d'ampleur dans les années à venir. Les challenges dans l'avenir seront de produire des modèles de plus en plus performants, capables de détecter au niveau de précision d'un agronome voire plus, les plantes agricoles et les différencier des mauvaises herbes puis de coupler ces modèles ultra performants à des robots chargés de détruire les mauvaises herbes.

Au regard de la diversité de plantes agricoles et tout autant de la grande variété de mauvaises herbes, une énorme quantité de données sera nécessaire pour créer une intelligence

artificielle “générale” capable de différencier les bonnes plantes des mauvaises ou mettre en place de plus petits modèles pour chaque région d’un pays ou chaque pays du monde.

## 7 Conclusion

Nous avons été heureux et enthousiastes tout au long de notre étude de travailler sur ce domaine porteur qu’est l’agriculture de précision. Nous tenons à réitérer nos remerciements auprès des chercheurs de AIT Austrian Institute of Technology qui ont réuni les datasets sur lesquels nous avons travaillé. Le dataset Crop0rWeed9 si bien organisé et annoté, fruit de l’OpenSource a permis à notre équipe de faire des avancées notables sur l’apprentissage profond particulières les réseaux neuronaux convolutionnels. Cette première étape de notre travail sera les prémices d’un travail futur plus large. Nous espérons voir d’autres équipes de chercheurs y contribuer de façon open source afin de permettre le développement de l’agriculture.

### *Remerciements :*

Nous aimerions remercier le Professeur Anass BELCAID, Professeur à l’École Nationale des Sciences Appliquées de Tétouan, pour ses différents conseils, sa pédagogie parfaite et ainsi que son aisance à partager les informations substantielles qui ont permis de réussir ce projet.

Code Source : [YOLOv3 Faster-RCNN](#)